

NATURE RESERVE SELECTION PROBLEM: A TIGHT
APPROXIMATION ALGORITHM

Magnus Bordewich and Charles Semple

*Department of Mathematics and Statistics
University of Canterbury
Private Bag 4800
Christchurch, New Zealand*

Report Number: UCDMS2007/1

FEBRUARY 2007

NATURE RESERVE SELECTION PROBLEM: A TIGHT APPROXIMATION ALGORITHM

MAGNUS BORDEWICH AND CHARLES SEMPLE

ABSTRACT. The Nature Reserve Selection Problem is a problem that arises in the context of studying biodiversity conservation. Subject to budgetary constraints, the problem is to select a set of regions to conserve so that the phylogenetic diversity of the set of species contained within those regions is maximized. Recently, it was shown by Moulton *et al.* [6] that this problem is NP-hard. In this paper, we establish a tight polynomial-time approximation algorithm for the Nature Reserve Selection Problem. Furthermore, we resolve a question on the computational complexity of a related problem left open in [6].

1. INTRODUCTION

Phylogenetic diversity is a quantitative tool for measuring the ‘biodiversity’ of a collection of species. This measure is based on the evolutionary distance amongst the species in the collection. Loosely speaking, if T is a phylogenetic tree whose leaves represent a set X of species and whose edges have real-valued lengths (weights), then the phylogenetic diversity (PD score) of a subset S of X is the sum of the weights of the edges of the minimal subtree of T connecting the species in S . The basic PD optimization problem is to find a subset of X of a given size that maximizes the PD score amongst all subsets of X of that size. Perhaps surprisingly, the greedy algorithm solves this problem exactly [1, 7, 11].

A natural extension of the basic problem allows for the consideration of conserving various regions such as nature reserves at some cost (see [6, 8, 9]). In particular, as well as an edge-weighted phylogenetic tree T , we have a collection \mathcal{A} of regions or areas containing species in X with each region having an associated cost of preservation. Given a fixed budget B , the PD optimization problem for this extension is to find a subset of the regions in \mathcal{A} to preserve that maximizes the PD score of the species contained within at least one preserved region while keeping within the budget. This problem is called Budgeted Nature Reserve Selection problem (BNRS).

Date: 5 February 2007.

1991 Mathematics Subject Classification. 05C05; 92D15.

Key words and phrases. Phylogenetic diversity; biodiversity conservation.

The first author was supported by the EPSRC, while the second author was supported by the New Zealand Marsden Fund.

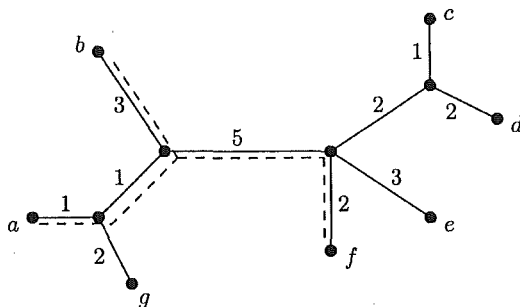


FIGURE 1. A phylogenetic X -tree with edge lengths, where $X = \{a, b, c, d, e, f, g\}$.

Moulton *et al.* [6] showed that this extension is NP-hard, that is, there is no polynomial-time algorithm for solving it unless $P=NP$. Despite this negative result, in this paper we show that there is a polynomial-time $(1 - 1/e)$ -approximation algorithm for this problem. That is, an efficient algorithm that generates a solution which has at least a $(1 - 1/e)$ fraction ($\approx 63\%$) of the phylogenetic diversity of the optimal solution. Moreover, this approximation ratio is the best possible.

The paper is arranged as follows. Section 2 contains a formal definition of BNRS and a discussion of related work. Section 3 contains the description of the approximation algorithm, and the statement of the main theorem, the proof of which is established in Section 4. Lastly, Section 5 answers a computational complexity question on a related problem that was left open in [6]. The notation and terminology in the paper follows [10].

2. BUDGETED NATURE RESERVE SELECTION

In order to define BNRS formally, we require the following definitions. A *phylogenetic X -tree* \mathcal{T} is an (unrooted) tree with no degree-2 vertices and whose leaf set is X . Let \mathcal{T} be a phylogenetic X -tree with edge set E and let $\lambda : E \rightarrow \mathbb{R}^{\geq 0}$ be an assignment of lengths (weights) to the edges of \mathcal{T} . Ignoring the dashed edges, Fig. 1 illustrates a phylogenetic X -tree with non-negative real-valued edge weights, where $X = \{a, b, c, d, e, f, g\}$.

For a subset S of X , the *phylogenetic diversity* (PD) of S on \mathcal{T} is the sum of the edge lengths of the minimal subtree of \mathcal{T} that connects S . This sum is denoted as $PD_{(\mathcal{T}, \lambda)}(S)$, however, if there is no ambiguity, we usually shorten it to $PD(S)$. Referring to Fig. 1, if $S = \{a, b, f\}$, then $PD(S)$ is equal to the sum of the weights of the minimal subtree (dashed edges) that connects a , b , and f , in particular, $PD(S) = 12$.

BNRS is formally defined as follows.

Problem: BNRS

Instance: A phylogenetic X -tree \mathcal{T} , a positive (real valued) weighting λ on the

edges of \mathcal{T} , a collection \mathcal{A} of subsets of X , a cost function c on the sets in \mathcal{A} , and a budget B .

Question: Find a subset \mathcal{A}' of \mathcal{A} that maximizes the PD score of $\bigcup_{A \in \mathcal{A}'} A$ on \mathcal{T} such that $\sum_{A \in \mathcal{A}'} c(A) \leq B$.

Referring to the informal discussions in the introduction, in the statement of BNRS, \mathcal{A} is the collection of regions and \mathcal{A}' is an optimal subset of regions that we wish to conserve that maximizes the PD score of the species contained in at least one of the preserved regions. Of course, the total cost of the preserving the regions in \mathcal{A}' is at most B .

The problem BNRS extends the problem OPTIMIZING DIVERSITY VIA REGIONS described in [6]. The extension from the latter to the former is that, instead of each region having a unit cost, the cost of conserving each region varies. Moulton *et al.* [6] showed that OPTIMIZING DIVERSITY VIA REGIONS is NP-hard and so, consequently, BNRS is also NP-hard. BNRS also extends the problem BUDGETED MAXIMUM COVERAGE, in which each element of X has a weight and the objective is to maximize the total weight of $\bigcup_{A \in \mathcal{A}'} A$ without the additional structure imposed by a tree [5]. An instance of the latter problem may be realized as a BNRS instance by taking \mathcal{T} to be a star tree with leaf set X and assigning the weight of each element in X to be the length of the incident edge in \mathcal{T} . (Note that a star tree is a phylogenetic tree with a single interior vertex.) The approximation algorithm and its proof presented here closely follow those in [5] for the restricted ‘star tree problem’, but must be extended to cover the more complicated interactions of PD score rather than a simple sum of weights. Lastly, BNRS is the “0 \leq 0/1 Nature Reserve Problem” briefly discussed in the appendix in [8].

3. THE APPROXIMATION ALGORITHM

In this section, we describe a tight polynomial-time approximation algorithm for BNRS called ApproxBNRS. The fact that it is such an algorithm is established in the next section. For a subset G of \mathcal{A} , the notations $c(G)$ and $PD(G)$ denote $\sum_{A \in G} c(A)$ and $PD(\bigcup_{A \in G} A)$, respectively.

We begin with an informal overview of ApproxBNRS and its subroutine Greedy (see Figs 2 and 3). By considering all possibilities, ApproxBNRS initially finds a feasible solution of size at most two that maximizes the PD score on \mathcal{T} . The resulting solution is called H_1 . Next, the algorithm, in turn, considers every subset of \mathcal{A} of size three and applies the subroutine Greedy to each of these subsets. The algorithm Greedy is a greedy-like algorithm that takes a subset G_0 of size three of \mathcal{A} and sequentially adds sets from $\mathcal{A} - G_0$. The only criteria for which set is selected is that, amongst all available sets, the ratio of incremental diversity to cost is maximized and we keep within budget. The resulting feasible solution that maximizes the PD score is called H_2 . Finally, ApproxBNRS compares the two feasible solutions H_1 and H_2 , and returns the one with the biggest PD score.

The main result of this paper is the following theorem whose proof is given in the next section.

```

Greedy( $G_0, U$ ):
 $G \leftarrow G_0$ 
Repeat
  select  $A \in U$  that maximizes  $\frac{PD(G \cup A) - PD(G)}{c(A)}$ 
  if  $c(G) + c(A) \leq B$  then
     $G \leftarrow G \cup \{A\}$ 
   $U \leftarrow U \setminus A$ 
Until  $U = \emptyset$ 
Return  $G$ 

```

FIGURE 2. The greedy algorithm Greedy.

```

ApproxBNRS( $T, \lambda, \mathcal{A}, c, B$ ):
Find  $G'$  in  $\{G : G \subseteq \mathcal{A}, c(G) \leq B, |G| \leq 2\}$  that maximizes PD
 $H_1 \leftarrow G'$ 
 $H_2 \leftarrow \emptyset$ 
For all  $G_0 \subseteq \mathcal{A}$ , such that  $|G_0| = 3$  and  $c(G_0) \leq B$  do
   $U \leftarrow \mathcal{A} \setminus G_0$ 
   $G \leftarrow \text{Greedy}(G_0, U)$ 
  if  $PD(G) > PD(H_2)$  then  $H_2 \leftarrow G$ 
If  $PD(H_1) > PD(H_2)$  then Return  $H_1$ , otherwise Return  $H_2$ 

```

FIGURE 3. The approximation algorithm ApproxBNRS.

Theorem 3.1. *ApproxBNRS is a polynomial-time $(1 - 1/e)$ -approximation algorithm for BNRS. Moreover, for any $\epsilon > 0$, BNRS cannot be approximated with an approximation ratio of $(1 - 1/e + \epsilon)$ unless $P=NP$.*

In terms of the running time of ApproxBNRS, running the greedy subroutine is very efficient, however, repeating this for all subsets of \mathcal{A} of size three incurs a multiplicative overhead of $O(|\mathcal{A}|^3)$. Typically the number of regions or nature reserves under consideration will be small, and hence this overhead is minor. Nevertheless, it is worth noting in the special case that all regions have the same cost, this term can be removed from the running time. In this situation, the greedy algorithm starting from a subset G_0 of \mathcal{A} of size two that maximizes the PD score amongst all 2-element subsets of \mathcal{A} achieves the approximation ratio $(1 - 1/e)$. The proof of this fact is a routine extension of [4], using the same insights regarding the difference between PD and the ordinary weight function as we have used in the proof of Theorem 3.1 given in the next section.

4. PROOF OF THEOREM 3.1

This section consists of the proof of Theorem 3.1. Let \mathcal{S}_{opt} denote the collection of subsets of \mathcal{A} of an optimal solution to BNRS. If $|\mathcal{S}_{\text{opt}}| \leq 2$, then ApproxBNRS finds

a feasible solution whose PD score is equal to the PD score of \mathcal{S}_{opt} . Therefore, we may assume that $|\mathcal{S}_{\text{opt}}| \geq 3$, in which case it suffices to show that there is a subset G_0 of \mathcal{A} with $|G_0| = 3$ whose input to Greedy (together with $\mathcal{A} - G_0$) results in a subset of \mathcal{A} whose PD score is within the approximation ratio stated in the theorem.

Let G_0 be the subset $\{S_1, S_2, S_3\}$ of \mathcal{S}_{opt} such that S_1 and S_2 are chosen to maximize $PD(S_1 \cup S_2)$ amongst all subsets of \mathcal{S}_{opt} of size two, and S_3 maximizes $PD(S_1 \cup S_2 \cup S_3)$ amongst all sets in $\mathcal{S}_{\text{opt}} \setminus \{S_1, S_2\}$. Now consider Greedy applied to $(G_0, \mathcal{A} - G_0)$. Let p denote the first iteration in which a member, A_{l+1} say, of $\mathcal{S}_{\text{opt}} - G_0$ is considered but, because of budgetary reasons, is not added to the current greedy solution. Up to iteration p , let, in order, A_1, A_2, \dots, A_l denote the members of $\mathcal{A} - G_0$ that are added to G_0 and, for $i = 1, \dots, l$, let $G_i = G_0 \cup \{A_1, A_2, \dots, A_i\}$. Observe that G_l is a feasible solution, and a subset of the final output G^* of the greedy subroutine, and hence $PD(G^*) \geq PD(G_l)$. For convenience, we also let $G_{l+1} = G_l \cup \{A_{l+1}\}$, but note that G_{l+1} is not a feasible solution as $c(G_{l+1}) > B$. Furthermore, for all i , let c_i denote $c(A_i)$. For a subset S of \mathcal{A} , denote the minimal subtree of \mathcal{T} that connects the elements of X that are contained in at least one member of S by $\mathcal{T}(S)$. We begin the proof with two lemmas.

Lemma 4.1. *For all $i \in \{1, 2, \dots, l+1\}$,*

$$PD(G_i) - PD(G_{i-1}) \geq \frac{c_i}{B} (PD(\mathcal{S}_{\text{opt}}) - PD(G_{i-1})).$$

Proof. One crucial point to observe in order for the approach of [5] to be applicable in our setting is that the incremental diversity from adding the entire optimal solution to the current partial greedy solution is bounded by the sum of the increments that would be obtained from adding each set in the optimal solution individually. We formalize this as follows. Let i be any element in $\{1, 2, \dots, l+1\}$. Let F denote the set of edges in $E(\mathcal{T}(\mathcal{S}_{\text{opt}} \cup G_{i-1})) - E(\mathcal{T}(G_{i-1}))$. Observe that $PD(\mathcal{S}_{\text{opt}} \cup G_{i-1}) - PD(G_{i-1})$ is equal to $\sum_{e \in F} \lambda(e)$. Since G_{i-1} is non-empty, there is, for each $e \in F$, an element in $\bigcup_{A \in (\mathcal{S}_{\text{opt}} - G_{i-1})} A$, such that e is on the path from that element to a vertex in $\mathcal{T}(G_{i-1})$. In particular, there is a set A_e in $\mathcal{S}_{\text{opt}} - G_{i-1}$ such that $\mathcal{T}(G_{i-1} \cup A_e)$ contains e . Since A_i is chosen so that $\frac{PD(G_i) - PD(G_{i-1})}{c_i}$ is maximized, we have, for all $A \in \mathcal{S}_{\text{opt}} - G_{i-1}$,

$$\frac{PD(G_{i-1} \cup A) - PD(G_{i-1})}{c(A)} \leq \frac{PD(G_i) - PD(G_{i-1})}{c_i}.$$

Therefore, as the total cost of the elements in $\mathcal{S}_{\text{opt}} - G_{i-1}$ is at most B ,

$$\begin{aligned}
PD(\mathcal{S}_{\text{opt}}) - PD(G_{i-1}) &\leq PD(\mathcal{S}_{\text{opt}} \cup G_{i-1}) - PD(G_{i-1}) \\
&= \sum_{e \in F} \lambda(e) \\
&\leq \sum_{A \in (\mathcal{S}_{\text{opt}} - G_{i-1})} \left[\sum_{\{e \in F: e \in T(G_{i-1} \cup A)\}} \lambda(e) \right] \\
&= \sum_{A \in (\mathcal{S}_{\text{opt}} - G_{i-1})} \frac{PD(G_{i-1} \cup A) - PD(G_{i-1})}{c(A)} c(A) \\
&\leq \sum_{A \in (\mathcal{S}_{\text{opt}} - G_{i-1})} \frac{PD(G_i) - PD(G_{i-1})}{c_i} c(A) \\
&\leq \frac{PD(G_i) - PD(G_{i-1})}{c_i} B.
\end{aligned}$$

Rearrangement now gives the inequality in the statement of the lemma and the result follows. \square

Lemma 4.2. *For all $i \in \{1, 2, \dots, l+1\}$,*

$$PD(G_i) - PD(G_0) \geq \left[1 - \prod_{k=1}^i \left(1 - \frac{c_k}{B} \right) \right] (PD(\mathcal{S}_{\text{opt}}) - PD(G_0)).$$

Proof. The proof is by induction on i . The result for $i = 1$ immediately follows from Lemma 4.1.

Now assume that $i \geq 2$ and that the result holds for all j , where $j < i$. Then, by Lemma 4.1 (for the first inequality) and induction (for the second inequality), we have

$$\begin{aligned}
PD(G_i) - PD(G_0) &= PD(G_{i-1}) - PD(G_0) + PD(G_i) - PD(G_{i-1}) \\
&\geq PD(G_{i-1}) - PD(G_0) + \frac{c_i}{B} (PD(\mathcal{S}_{\text{opt}}) - PD(G_{i-1})) \\
&= PD(G_{i-1}) - PD(G_0) \\
&\quad + \frac{c_i}{B} (PD(\mathcal{S}_{\text{opt}}) - PD(G_0) - (PD(G_{i-1}) - PD(G_0))) \\
&= \left(1 - \frac{c_i}{B} \right) (PD(G_{i-1}) - PD(G_0)) + \frac{c_i}{B} (PD(\mathcal{S}_{\text{opt}}) - PD(G_0)) \\
&\geq \left(1 - \frac{c_i}{B} \right) \left[1 - \prod_{k=1}^{i-1} \left(1 - \frac{c_k}{B} \right) \right] (PD(\mathcal{S}_{\text{opt}}) - PD(G_0)) \\
&\quad + \frac{c_i}{B} (PD(\mathcal{S}_{\text{opt}}) - PD(G_0)) \\
&= \left[1 - \prod_{k=1}^i \left(1 - \frac{c_k}{B} \right) \right] (PD(\mathcal{S}_{\text{opt}}) - PD(G_0)).
\end{aligned}$$

This completes the proof of the lemma. \square

Before completing the proof of Theorem 3.1, we define the problem MAXIMUM k -COVERAGE which we refer to in establishing the second part of the theorem.

Problem: MAXIMUM k -COVERAGE

Instance: A collection \mathcal{A} of subsets of X and an integer k .

Question: Find a subset $\mathcal{A}' = \{A_1, A_2, \dots, A_k\}$ of \mathcal{A} of size k that maximizes the size of the set $A_1 \cup A_2 \cup \dots \cup A_k$.

Feige [2] showed that no polynomial-time approximation algorithm for MAXIMUM k -COVERAGE can have an approximation ratio better than $(1 - 1/e)$ unless $P=NP$.

Proof of Theorem 3.1. Since $c(G_{l+1}) > B$ and the function $\prod_{k=1}^{l+1} \left(1 - \frac{c_k}{c(G_{l+1})}\right)$, where $c(G_{l+1}) = \sum_k c_k$, has a maximum at $c_k = c(G_{l+1})/(l+1)$ for all k , we have

$$\begin{aligned} 1 - \prod_{k=1}^{l+1} \left(1 - \frac{c_k}{B}\right) &\geq 1 - \prod_{k=1}^{l+1} \left(1 - \frac{c_k}{c(G_{l+1})}\right) \\ &\geq 1 - \left(1 - \frac{1}{l+1}\right)^{l+1} \\ &\geq 1 - 1/e. \end{aligned}$$

Hence, by Lemma 4.2, we have

$$(1) \quad PD(G_{l+1}) - PD(G_0) \geq (1 - 1/e)(PD(S_{\text{opt}}) - PD(G_0)).$$

Recalling that $G_0 = \{S_1, S_2, S_3\}$, we now show that

$$(2) \quad PD(S_1 \cup S_2 \cup S_3) - PD(S_1 \cup S_2) \leq PD(G_0)/3.$$

Let $A_j = E(T(S_1 \cup S_2 \cup S_3)) - E(T((S_1 \cup S_2 \cup S_3) - S_j))$ for $j = 1, 2, 3$. Since

$$\begin{aligned} PD(S_1 \cup S_2 \cup S_3) &= PD(S_1 \cup S_2) + \sum_{e \in A_3} \lambda(e) \\ &= PD(S_1 \cup S_3) + \sum_{e \in A_2} \lambda(e) \\ &= PD(S_2 \cup S_3) + \sum_{e \in A_1} \lambda(e), \end{aligned}$$

and since S_1 and S_2 were chosen to maximize $PD(S_1 \cup S_2)$, it follows that

$$\sum_{e \in A_3} \lambda(e) \leq \sum_{e \in A_j} \lambda(e) \quad j = 1, 2.$$

It is easily seen that each edge in $E(T(S_1 \cup S_2 \cup S_3))$ occurs in at most one A_j . Hence

$$\begin{aligned} PD(S_1 \cup S_2 \cup S_3) &\geq \sum_{j=1}^3 \sum_{e \in A_j} \lambda(e) \\ &\geq 3 \sum_{e \in A_3} \lambda(e), \end{aligned}$$

and so

$$PD(S_1 \cup S_2 \cup S_3) - PD(S_1 \cup S_2) = \sum_{e \in A_3} \lambda(e) \leq PD(G_0)/3,$$

giving Eqn (2).

Next,

$$PD(G_{l+1}) - PD(G_l) \leq PD(S_1 \cup S_2 \cup A_{l+1}) - PD(S_1 \cup S_2),$$

and so

$$(3) \quad PD(G_{l+1}) - PD(G_l) \leq PD(S_1 \cup S_2 \cup S_3) - PD(S_1 \cup S_2) \leq PD(G_0)/3;$$

otherwise A_{l+1} would have been chosen instead of S_3 to be in G_0 . Putting together Eqns (1) and (3), we get

$$\begin{aligned} PD(G_l) &\geq PD(G_{l+1}) - PD(G_0)/3 \\ &\geq (1 - 1/e)(PD(\mathcal{S}_{\text{opt}}) - PD(G_0)) + (1 - \frac{1}{3})PD(G_0) \\ &> (1 - 1/e)PD(\mathcal{S}_{\text{opt}}). \end{aligned}$$

This proves the first part of the theorem.

For the proof of the second part, we observe that BNRS is a generalization of MAXIMUM k -COVERAGE and so, by Feige's result, no approximation algorithm can exist for BNRS with ratio better than $(1 - 1/e)$ unless $P=NP$.

Given an instance of MAXIMUM k -COVERAGE, take T to be the star tree on leaf set X in which each edge has weight 1. Assign a cost of 1 to each element of \mathcal{A} and take the budget $B = k$. Under this set-up, it is clear that MAXIMUM k -COVERAGE can be interpreted as a special case of BNRS. Hence a polynomial-time approximation algorithm for BNRS with approximation ratio α would yield an approximation algorithm for MAXIMUM k -COVERAGE with approximation ratio α . By [2], no such algorithm can exist for $\alpha = (1 - 1/e + \epsilon)$ unless $P=NP$. \square

5. OPTIMIZING DIVERSITY WITH COVERAGE

The problem OPTIMIZING DIVERSITY WITH COVERAGE was defined in [6], where a very restricted version was shown to have a polynomial-time algorithm. While superficially this problem is similar to BNRS, the problem behaves very differently. Loosely speaking, we are given an edge-weighted phylogenetic X -tree T and a collection \mathcal{A} of subsets of X . Here the members of \mathcal{A} represent some attributes that the species possess. For example, $\mathcal{A} = \{A_1, A_2, \dots, A_s\}$ may be a collection of taxonomic groups and each A_i contains the species in X that belong to the group. Given a fixed positive integer k and positive integers n_1, n_2, \dots, n_s , the PD optimization problem is to find a subset X' of X of size k that contains, for all i , at least n_i species with attribute A_i and maximizes the PD score amongst all such subsets of X of size k . Formally, we have the following problem.

Problem: OPTIMIZING DIVERSITY WITH COVERAGE

Instance: A phylogenetic X -tree T , a positive real-valued weighting λ on the edges of T , a collection \mathcal{A} of subsets of X , a threshold n_A for each $A \in \mathcal{A}$, and a positive

integer k .

Question: Find a subset X' of X that maximizes the PD score of X' on \mathcal{T} such that $|X'| \leq k$ and, for each $A \in \mathcal{A}$, at least n_A species from A are included in X' .

The restricted case solved in [6] is when each element of X appears in exactly one set $A \in \mathcal{A}$ and the subtrees in $\{\mathcal{T}(A) : A \in \mathcal{A}\}$ are vertex disjoint. While this restricted version is shown to be solvable in polynomial time, the question of the computational complexity of the problem under less stringent or no restrictions is left open. We end this section by observing that determining if there is even a feasible solution to the general problem OPTIMIZING DIVERSITY WITH COVERAGE is NP-hard, let alone finding an optimal solution. This is because determining if there is a feasible solution is equivalent to the classic NP-complete decision problem HITTING SET [3].

Problem: HITTING SET

Instance: A collection \mathcal{A} of subsets of X and an integer k .

Question: Does there exist a subset X' of X of size at most k such that $A \cap X' \neq \emptyset$ for all $A \in \mathcal{A}$.

For an instance of HITTING SET as above, consider the instance of OPTIMIZING DIVERSITY WITH COVERAGE by taking the same sets X and \mathcal{A} , and integer k . Now take $n_A = 1$ for all $A \in \mathcal{A}$ and let \mathcal{T} be an arbitrary phylogenetic X -tree. Then a subset of X is a feasible solution to the latter problem if and only if it is a feasible solution to the former problem. Conversely, for an instance of OPTIMIZING DIVERSITY WITH COVERAGE, consider the instance of HITTING SET by taking the ground set to be X , the bound to be k , and choosing the collection of subsets of X to be

$$\{B : \exists A \in \mathcal{A}, B \subseteq A, |B| = |A| - n_A + 1\}.$$

In words, this collection consists of, for each $A \in \mathcal{A}$, all subsets of A of size $|B| = |A| - n_A + 1$. It is now easily seen that a subset of X is a feasible solution to this instance of HITTING SET if and only if it is a feasible solution to the original instance of OPTIMIZING DIVERSITY WITH COVERAGE.

The above equivalence suggests that the restrictions required to make OPTIMIZING DIVERSITY WITH COVERAGE solvable, or even approximable, must be fairly severe. Certainly they must at least make the associated restricted version of HITTING SET tractable. One example could be to restrict k to be at least $\sum_{A \in \mathcal{A}} n_A$. In this case, HITTING SET is trivial, and hence a feasible solution to OPTIMIZING DIVERSITY WITH COVERAGE can be found easily. However, it is still not clear whether the optimal solution can be found efficiently.

REFERENCES

- [1] Faith, D.P., 1992. Conservation evaluation and phylogenetic diversity. *Biol. Conserv.* 61, 1-10.
- [2] Feige, U., 1998. A threshold of $\ln n$ for approximating set cover. *J. ACM* 45, 634-652.
- [3] Garey, M.R., Johnson, D.S., 1979. *Computers and intractability: A guide to the theory of NP-completeness*, Freeman, San Francisco, CA.
- [4] Hochbaum, D., 1997. Approximating covering and packing problems: set cover, vertex cover, independent set, and related problems. In: Hochbaum, D. (Ed.), *Approximation Algorithms for NP-Hard Problems*. PWS, Boston.

- [5] Khuller, S., Moss, A., Naor, J., 1999. The budgeted maximum coverage problem. Inform. Process. Lett. 70, 39-45.
- [6] Moulton, V., Semple, C., Steel, M., 2007. Optimizing phylogenetic diversity under constraints. J. Theoret. Biol., in press.
- [7] Pardi, F., Goldmann, N., 2005. Species choice for comparative genomics: being greedy works. PLoS Genetics 1, e71.
- [8] Pardi, F., Goldman, N., 2007. Resource aware taxon selection for maximising phylogenetic diversity. Syst. Biol., in press.
- [9] Rodrigues, A.S.L., Brooks, T.M., Gaston, K.J., 2005. Integrating phylogenetic diversity in the selection of priority areas for conservation: does it make a difference In Purvis, A., Gittleman, J.L., Brooks, T. (Eds.), Phylogeny and Conservation. Cambridge University Press, Cambridge.
- [10] Semple, C., Steel, M., 2003. Phylogenetics. Oxford University Press, Oxford.
- [11] Steel, M., 2005. Phylogenetic diversity and the greedy algorithm. Syst. Biol. 54, 527-529.

DEPARTMENT OF COMPUTER SCIENCE, DURHAM UNIVERSITY, DURHAM DH1 3LE, UNITED KINGDOM

E-mail address: m.j.r.bordewich@durham.ac.uk

BIOMATHEMATICS RESEARCH CENTRE, DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF CANTERBURY, CHRISTCHURCH, NEW ZEALAND

E-mail address: c.semple@math.canterbury.ac.nz